

日 本 国 特 許 庁  
JAPAN PATENT OFFICE

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日            2 0 0 3 年   2 月 1 9 日  
Date of Application:

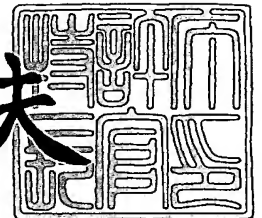
出 願 番 号            特 願 2 0 0 3 - 0 4 1 4 8 6  
Application Number:  
[ST. 10/C]:            [ J P 2 0 0 3 - 0 4 1 4 8 6 ]

出   願   人            株 式 会 社 東 芝  
Applicant(s):

2 0 0 3 年   7 月 1 8 日

特許庁長官  
Commissioner,  
Japan Patent Office

今 井 康 夫



【書類名】 特許願

【整理番号】 A000205441

【提出日】 平成15年 2月19日

【あて先】 特許庁長官 殿

【国際特許分類】 G06F 15/00

【発明の名称】 ストレージ装置、分担範囲決定方法及びプログラム

【請求項の数】 20

【発明者】

    【住所又は居所】 神奈川県川崎市幸区小向東芝町 1 番地 株式会社東芝研  
    究開発センター内

    【氏名】 吉田 英樹

【発明者】

    【住所又は居所】 神奈川県川崎市幸区小向東芝町 1 番地 株式会社東芝研  
    究開発センター内

    【氏名】 金井 達徳

【発明者】

    【住所又は居所】 神奈川県川崎市幸区小向東芝町 1 番地 株式会社東芝研  
    究開発センター内

    【氏名】 崎山 伸夫

【特許出願人】

    【識別番号】 000003078

    【氏名又は名称】 株式会社 東芝

【代理人】

    【識別番号】 100058479

    【弁理士】

    【氏名又は名称】 鈴江 武彦

    【電話番号】 03-3502-3181

## 【選任した代理人】

【識別番号】 100091351

【弁理士】

【氏名又は名称】 河野 哲

## 【選任した代理人】

【識別番号】 100088683

【弁理士】

【氏名又は名称】 中村 誠

## 【選任した代理人】

【識別番号】 100108855

【弁理士】

【氏名又は名称】 蔵田 昌俊

## 【選任した代理人】

【識別番号】 100084618

【弁理士】

【氏名又は名称】 村松 貞男

## 【選任した代理人】

【識別番号】 100092196

【弁理士】

【氏名又は名称】 橋本 良郎

## 【手数料の表示】

【予納台帳番号】 011567

【納付金額】 21,000円

## 【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【ブルーフの要否】 要

【書類名】 明細書

【発明の名称】 ストレージ装置、分担範囲決定方法及びプログラム

【特許請求の範囲】

【請求項 1】

分散ストレージシステムの構成に用いられるストレージ装置であって、  
前記分散ストレージシステムが対象とする識別子空間のうちの分担範囲に含まれる識別子を持つデータを格納するための手段と、  
前記識別子空間に割り当てられた、前記分担範囲を決定する基準となる基準位置を記憶する第 1 の記憶手段と、  
前記分担範囲を決定する際に考慮される重みを記憶する第 2 の記憶手段と、  
前記第 2 の記憶手段に記憶された重みと、前記識別子空間において前記分担範囲に隣接する範囲を分担する隣接ストレージ装置に係る重みとに基づいて、該隣接ストレージ装置とで前記識別子空間を分担する場合における相対的な空間幅を決定する第 1 の決定手段と、  
前記第 1 の決定手段により決定された前記相対的な空間幅に基づいて、前記識別子空間における前記第 1 の記憶手段に記憶された基準位置から前記隣接ストレージ装置に係る基準位置までの範囲内について分担する分担範囲を決定する第 2 の決定手段とを備えたことを特徴とするストレージ装置。

【請求項 2】

前記第 1 の決定手段は、前記第 2 の記憶手段に記憶された重みと、前記隣接ストレージ装置に係る重みとを加算し、これによって得られた値で前記第 2 の記憶手段に記憶された重みを除した値を、前記相対的な空間幅とすることを特徴とする請求項 1 に記載のストレージ装置。

【請求項 3】

前記第 2 の決定手段は、前記第 1 の記憶手段に記憶された前記基準位置を示す値に対して前記相対的な空間幅を乗じて得た値と、前記隣接ストレージ装置に係る前記基準位置を示す値に対して前記相対的な空間幅を 1 から減じた値を乗じて得た値とを、加算し、これによって得られた値が示す位置と、前記第 1 の記憶手段に記憶された前記基準位置との間の空間を、前記基準位置間の範囲内について

分担する分担範囲とすることを特徴とする請求項 2 に記載のストレージ装置。

**【請求項 4】**

前記隣接ストレージ装置として、第 1 の隣接ストレージ装置と第 2 の隣接ストレージ装置の 2 つが存在する場合には、該第 1 の隣接ストレージ装置について前記第 1 の決定手段による前記決定及び前記第 2 の決定手段による前記決定を行って、第 1 の分担範囲を決定するとともに、該第 2 の隣接ストレージ装置について前記第 1 の決定手段による前記決定及び前記第 2 の決定手段による前記決定を行って、第 2 の分担範囲を決定し、該第 1 の分担範囲と該第 2 の分担範囲とを包含する範囲を、求めるべき全分担範囲とすることを特徴とする請求項 1 に記載のストレージ装置。

**【請求項 5】**

前記隣接ストレージ装置として、唯一の隣接ストレージ装置のみが存在する場合には、該隣接ストレージ装置について前記第 1 の決定手段による前記決定及び前記第 2 の決定手段による前記決定を行って、第 1 の分担範囲を決定するとともに、前記第 1 の記憶手段に記憶された前記基準位置と、前記識別子空間の両端部のうち前記第 1 の記憶手段に記憶された前記基準位置に近い方の端部の位置との間の空間を、第 2 の分担範囲として決定し、該第 1 の分担範囲と該第 2 の分担範囲とを包含する範囲を、求めるべき全分担範囲とすることを特徴とする請求項 1 に記載のストレージ装置。

**【請求項 6】**

前記分散ストレージシステムを構成する他のストレージ装置との通信に用いるために割り当てられたアドレスを記憶する第 3 の記憶手段を更に備え、

前記アドレスに対して所定のハッシュ関数を適用して得られる値を、前記基準位置を示す値とすることを特徴とする請求項 1 に記載のストレージ装置。

**【請求項 7】**

前記隣接ストレージ装置に係るアドレスに対して前記所定のハッシュ関数を適用して得られる値を、前記隣接ストレージ装置に係る前記基準位置を示す値とすることを特徴とする請求項 6 に記載のストレージ装置。

**【請求項 8】**

前記分散ストレージシステムを構成する他の全部又は一部のストレージ装置に割り当てられたアドレスを示す情報を取得するための取得手段と、

前記取得手段により取得された前記他のストレージ装置のアドレスを記憶する第 4 の記憶手段とを更に備えたことを特徴とする請求項 6 または 7 に記載のストレージ装置。

**【請求項 9】**

前記隣接ストレージ装置に係るアドレス、及び前記第 3 の記憶手段に記憶された前記アドレスに対して所定のハッシュ関数を適用して得られる値に、1, 2, 4, 8, ...,  $2^{(b-1)}$  ( $b$  は予め定められた整数) をそれぞれ加えた値に相当する識別子を担当する前記他のストレージ装置にそれぞれ割り当てられた前記アドレスを、前記取得手段により予め取得し、前記第 4 の記憶手段に記憶しておくことを特徴とする請求項 8 に記載のストレージ装置。

**【請求項 10】**

前記アドレスは、IP アドレスであることを特徴とする請求項 6 ないし 9 のいずれか 1 項に記載のストレージ装置。

**【請求項 11】**

1 台の前記ストレージ装置を複数の仮想ノードに対応させ、各仮想ノードに異なる分担範囲を分担させるものとして、個々の仮想ノードの分担範囲を、前記第 1 及び第 2 の決定手段によりそれぞれ決定することを特徴とする請求項 1 ないし 10 のいずれか 1 項記載のストレージ装置。

**【請求項 12】**

前記ストレージ装置のアドレスに当該仮想ノードの識別情報を多重化し、これに対して所定のハッシュ関数を適用して得られる値を、当該仮想ノードの前記基準位置を示す値とすることを特徴とする請求項 11 に記載のストレージ装置。

**【請求項 13】**

各仮想ノードごとに重みを付与することを特徴とする請求項 11 または 12 に記載のストレージ装置。

**【請求項 14】**

全仮想ノードについて共通の重みを付与することを特徴とする請求項 11 また

は 12 に記載のストレージ装置。

【請求項 15】

前記識別子空間において前記第 1 の記憶手段に記憶された基準位置に対して最も近い基準位置を持つ他のストレージ装置を、前記隣接ストレージ装置とすることを特徴とする請求項 1 ないし 14 のいずれか 1 項に記載のストレージ装置。

【請求項 16】

前記識別子空間において前記第 1 の記憶手段に記憶された基準位置に対して n 番目 (n は予め定められた 2 以上の整数) に近い基準位置を持つ他のストレージ装置を、前記隣接ストレージ装置とみなして、前記第 1 の決定手段による前記決定及び前記第 2 の決定手段による前記決定を行うことを特徴とする請求項 1 ないし 14 のいずれか 1 項に記載のストレージ装置。

【請求項 17】

前記データは、ファイル又はファイルのブロックであることを特徴とする請求項 1 ないし 15 のいずれか 1 項に記載のストレージ装置。

【請求項 18】

分散ストレージシステムの構成に用いられるストレージ装置であって、

前記分散ストレージシステムが対象とする識別子空間のうちで自装置が分担する分担範囲に含まれる識別子を持つデータを格納するための手段と、

自装置に付与された重みと、前記識別子空間における自装置の分担範囲に隣接する範囲を分担する隣接ストレージ装置に付与された重みに基づいて、自装置と該隣接ストレージ装置とで前記識別子空間を分担する場合において自装置が分担する相対的な空間幅を決定する第 1 の決定手段と、

前記第 1 の決定手段により決定された前記相対的な空間幅に基づいて、自装置に予め割り当てられた前記識別子空間における基準位置から前記隣接ストレージ装置に予め割り当てられた前記識別子空間における基準位置までの範囲内について自装置が分担する分担範囲を決定する第 2 の決定手段とを備えたことを特徴とするストレージ装置。

【請求項 19】

分散ストレージシステムの構成に用いられ、該分散ストレージシステムが対象

とする識別子空間のうちの分担範囲に含まれる識別子を持つデータを格納するストレージ装置における分担範囲決定方法であって、

前記識別子空間に割り当てられた、前記分担範囲を決定する基準となる基準位置を、第1の記憶手段に記憶するステップと、

前記分担範囲を決定する際に考慮される重みを、第2の記憶手段に記憶するステップと、

前記第2の記憶手段に記憶された重みと、前記識別子空間において前記分担範囲に隣接する範囲を分担する隣接ストレージ装置に係る重みに基づいて、該隣接ストレージ装置とで前記識別子空間を分担する場合における相対的な空間幅を決定する第1の決定ステップと、

前記第1のステップ手段により決定された前記相対的な空間幅に基づいて、前記識別子空間における前記第1の記憶手段に記憶された基準位置から前記隣接ストレージ装置に係る基準位置までの範囲内について分担する分担範囲を決定する第2の決定ステップとを有することを特徴とする分担範囲決定方法。

#### 【請求項20】

分散ストレージシステムの構成に用いられるストレージ装置としてコンピュータを機能させるためのプログラムであって、

前記分散ストレージシステムが対象とする識別子空間のうちの分担範囲に含まれる識別子を持つデータを格納するための機能と、

前記識別子空間に割り当てられた、前記分担範囲を決定する基準となる基準位置を記憶する第1の記憶機能と、

前記分担範囲を決定する際に考慮される重みを記憶する第2の記憶機能と、

前記第2の記憶機能に記憶された重みと、前記識別子空間において前記分担範囲に隣接する範囲を分担する隣接ストレージ装置に係る重みに基づいて、該隣接ストレージ装置とで前記識別子空間を分担する場合における相対的な空間幅を決定する第1の決定機能と、

前記第1の決定機能により決定された前記相対的な空間幅に基づいて、前記識別子空間における前記第1の記憶機能に記憶された基準位置から前記隣接ストレージ装置に係る基準位置までの範囲内について分担する分担範囲を決定する第2



の決定機能とを実現させるためのプログラム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、分散ストレージシステムの構成に用いられるストレージ装置、分散ストレージシステムが対象とする識別子空間のうちで自装置が分担する分担範囲を決定するための分担範囲決定方法及びプログラムに関する。

【0002】

【従来の技術】

近年、ネットワーク上に分散配置された計算機で処理を分担するグリッドコンピューティングや、災害への対応などへの要求から、広域ネットワークにストレージを多数分散配置して一つの仮想的なストレージを実現する、という分散ストレージシステムが注目されている。

【0003】

このようなシステムでは、システムを構成する各ストレージノードが増設や故障などによって頻繁に増減するため、ファイルシステムまたはファイルを手動で個別にストレージに割り当てるのは非現実的である。また、集中制御を行なうサーバを設けてファイルのストレージノードへの割り当てを集中管理する方式には、そのサーバの障害発生や負荷集中がシステム全体に影響するという問題がある。そのため、ファイルのストレージノードへの自動的な割り当てを、集中制御なしにどうやって分散して行なうかが問題とされていた。

【0004】

この問題を解決するため、非特許文献1や非特許文献2のような研究的な分散ストレージシステムでは、ストレージノードのアドレスにハッシュ関数を適用してノードIDを決めることによって、ストレージノードをファイルIDと同じ空間にマッピングし、当該ファイルの持つファイルIDに最も近いノードIDを持つストレージノードに、当該ファイルを割り当てている。

【0005】

この方式では、ファイルをどのストレージノードに格納するかは、他のストレ

ージノードのノードIDのリストさえあれば、一つのストレージノード内での計算で決定できる。そのため、ファイルの割り当てを集中制御するサーバが不要になるほか、ファイルの割り当てをファイル一つ一つについてストレージノード間で個別に調整する必要もなくなる。よって、ストレージノードの増減の際にのみストレージノードのアドレスを他のストレージノードへ通知すればよく、ストレージノード間の通信量が減るとともに処理の並列性が向上する。

#### 【0006】

ストレージノードが追加された場合には、追加されたストレージノードとハッシュ空間上で隣接するストレージノードとの中間点で空間の再分割を行ない、以前からあったストレージノードの担当していた空間の一部を新しいストレージノードが引き継いで担当する。逆にストレージノードがなくなる場合には、なくなるストレージノードが担当している空間を分割して両隣のストレージノードに追加分担してもらう。

#### 【0007】

特定のストレージノードにファイルが集中しないよう、ファイルIDはファイルID空間上に一様に分散させる必要があるため、ハッシュ関数を用いてファイルIDをつけるのが一般的である。ハッシュ関数の引数としては、ファイル名を使う場合や、ファイルの中身のデータを使う場合などがある。

#### 【0008】

ファイルは、CFSのようにブロック単位に分割してブロックごとにストレージノードに配置してもよいし、CANのようにファイル全体をまとめてストレージノードに配置してもよい。

#### 【0009】

##### 【非特許文献1】

CFS (Wide-area cooperative storage with CFS, Frank Dabek, M. Frans Kaashoek, David Karger, Robert Morris, and Ion Stoica, 18th ACM Symposium on Operating Systems Principles (SOSP '01), October 2001)

#### 【0010】

##### 【非特許文献2】

CAN (A Scalable Content-Addressable Network, Sylvia Ratnasamy, Paul Francis, Mark Handley, Richard Karp and Scott Shenker, ACM SIGCOMM 2001)

【0011】

【発明が解決しようとする課題】

しかし、このようなハッシュによるファイルのストレージノードへの割り当てでは、個々のストレージノードに割り当てられる空間の大きさの期待値が全ストレージノードで均等となる。このため、記憶容量・演算能力・回線速度などがストレージノードごとに異なる場合に対応できない。これによって、容量の大きなストレージノードで容量が余っているのに、容量の小さなストレージノードでの容量不足が生じ、システム全体としてファイルの保存ができなくなったり、演算能力の低いストレージノードにも能力の高いストレージノードと同じ量のI/Oが要求されて応答速度が低下したりする、という問題が生じる。分散ストレージシステムを研究レベルを超えて実用に供する際には、この点が障害となる。

【0012】

CFSでは、この問題を回避するため、容量の大きなストレージノードを、仮想ノードと呼ばれる複数のストレージノードに仮想的に対応させることを提案している。確かに、容量の大きなストレージノードが容量の小さなストレージノードのたとえば数倍程度であれば、大きなストレージノードを数個に分割すればよい。しかし、ストレージノード間で容量が大きく異なる場合、例えば、容量の大きいストレージノードが容量の小さいストレージノードの数千倍の容量を持つ場合は、容量の大きなストレージノードを数千個の仮想ノードに分割する必要がある、仮想ノードの管理オーバーヘッドが問題となる。また、ディスク技術の進歩などによってストレージノードの平均的な容量が変化した場合に、仮想ノードの単位をどう調整するかも問題となる。したがって、仮想ノード単体では容量の多様性への対応が不十分である。

【0013】

なお、ここでは、ファイル全体をまとめてストレージノードに配置した分散ストレージを中心に述べたが、ファイルをブロック単位に分割してブロックごとにストレージノードに配置した分散ストレージについても同様の問題がある（ファ



イル名とブロック番号との組をブロックの名前として管理する場合を考えれば同様である)。

#### 【0014】

本発明は、上記事情を考慮してなされたもので、分散ストレージシステムを構成する各ストレージ装置が、該分散ストレージシステムが対象とする識別子空間のうちで分担する分担範囲を決定するにあたって、より効果的な分担範囲の決定を行うことができるストレージ装置、分担範囲決定方法及びプログラムを提供することを目的とする。

#### 【0015】

##### 【課題を解決するための手段】

本発明は、分散ストレージシステムの構成に用いられるストレージ装置であって、前記分散ストレージシステムが対象とする識別子空間のうちの分担範囲に含まれる識別子を持つデータを格納するための手段と、前記識別子空間に割り当てられた、前記分担範囲を決定する基準となる基準位置を記憶する第1の記憶手段と、前記分担範囲を決定する際に考慮される重みを記憶する第2の記憶手段と、前記第2の記憶手段に記憶された重みと、前記識別子空間において前記分担範囲に隣接する範囲を分担する隣接ストレージ装置に係る重みとに基づいて、該隣接ストレージ装置とで前記識別子空間を分担する場合における相対的な空間幅を決定する第1の決定手段と、前記第1の決定手段により決定された前記相対的な空間幅に基づいて、前記識別子空間における前記第1の記憶手段に記憶された基準位置から前記隣接ストレージ装置に係る基準位置までの範囲内について分担する分担範囲を決定する第2の決定手段とを備えたことを特徴とする。

#### 【0016】

また、本発明は、分散ストレージシステムの構成に用いられるストレージ装置であって、前記分散ストレージシステムが対象とする識別子空間のうちで自装置が分担する分担範囲に含まれる識別子を持つデータを格納するための手段と、自装置に付与された重みと、前記識別子空間における自装置の分担範囲に隣接する範囲を分担する隣接ストレージ装置に付与された重みとに基づいて、自装置と該隣接ストレージ装置とで前記識別子空間を分担する場合において自装置が分担す

る相対的な空間幅を決定する第1の決定手段と、前記第1の決定手段により決定された前記相対的な空間幅に基づいて、自装置に予め割り当てられた前記識別子空間における基準位置から前記隣接ストレージ装置に予め割り当てられた前記識別子空間における基準位置までの範囲内について自装置が分担する分担範囲を決定する第2の決定手段とを備えたことを特徴とする。

#### 【0017】

本発明では、記憶容量・演算能力・回線速度などをもとにストレージ装置の重みを決め、その重みに応じた空間割り当てを行なう。例えば、自装置と隣接ストレージ装置との重みの比に比例した幅になるように、自ノードと隣接ストレージ装置との間の空間を二つに分割する。これによって、重みに比例等した幅の空間をストレージ装置に割り当てることができる。

#### 【0018】

好ましくは、1台の前記ストレージ装置を複数の仮想ノードに対応させ、各仮想ノードに相異なる分担範囲を分担させるものとして、個々の仮想ノードの分担範囲を、前記第1及び第2の決定手段によりそれぞれ決定するようにしてもよい。このように、一つの物理ノードに対して複数の仮想ノードを割り当てることによって、ストレージ装置への割り当て空間をより正確に重みに比例させることができる。仮想化を行なえば、一つの物理ノードに対応するすべての仮想ノードについて同じような重みのストレージ装置とばかり隣接する確率は低いため、一方の仮想ノードについて広い幅を割り当てられても、幅の合計は重みに比例した大きさとなる（重みが小さなストレージ装置ばかりが隣接した空間に集中して容量が小さいにも関わらず大きな幅の空間を割り当てられてしまうようなことが回避できる）。

#### 【0019】

また、好ましくは、前記識別子空間において前記第1の記憶手段に記憶された基準位置に対してn番目（nは予め定められた2以上の整数）に近い基準位置を持つ他のストレージ装置を、前記隣接ストレージ装置とみなして、前記第1の決定手段による前記決定及び前記第2の決定手段による前記決定を行うようにしてもよい。分散ストレージの目的の一つである信頼性を実現するために、一つのデ

ータを複数のストレージ装置に置くという冗長化を行なう場合があるが、このようにすることによって、重みを考慮した冗長化を行なう場合が可能となる。一つのデータを  $n$  個のストレージノードで多重化する場合には、基準位置で定義される隣のストレージ装置とではなく、 $n$  個先のストレージ装置との間で重みに応じた空間分割を行ない、一つの点が  $n$  個のストレージ装置に重複してマッピングされるようにする。これによって、一部のストレージ装置が故障・災害・回線障害によって使えなくなった場合でも、同じ空間を担当する他のストレージ装置を利用してデータを読み書きすることができる。

#### 【0020】

また、本発明は、分散ストレージシステムの構成に用いられ、該分散ストレージシステムが対象とする識別子空間のうちの分担範囲に含まれる識別子を持つデータを格納するストレージ装置における分担範囲決定方法であって、前記識別子空間に割り当てられた、前記分担範囲を決定する基準となる基準位置を、第1の記憶手段に記憶するステップと、前記分担範囲を決定する際に考慮される重みを、第2の記憶手段に記憶するステップと、前記第2の記憶手段に記憶された重みと、前記識別子空間において前記分担範囲に隣接する範囲を分担する隣接ストレージ装置に係る重みとに基づいて、該隣接ストレージ装置とで前記識別子空間を分担する場合における相対的な空間幅を決定する第1の決定ステップと、前記第1のステップ手段により決定された前記相対的な空間幅に基づいて、前記識別子空間における前記第1の記憶手段に記憶された基準位置から前記隣接ストレージ装置に係る基準位置までの範囲内について分担する分担範囲を決定する第2の決定ステップとを有することを特徴とする。

#### 【0021】

また、本発明は、分散ストレージシステムの構成に用いられるストレージ装置としてコンピュータを機能させるためのプログラムであって、前記分散ストレージシステムが対象とする識別子空間のうちの分担範囲に含まれる識別子を持つデータを格納するための機能と、前記識別子空間に割り当てられた、前記分担範囲を決定する基準となる基準位置を記憶する第1の記憶機能と、前記分担範囲を決定する際に考慮される重みを記憶する第2の記憶機能と、前記第2の記憶機能に

記憶された重みと、前記識別子空間において前記分担範囲に隣接する範囲を分担する隣接ストレージ装置に係る重みに基づいて、該隣接ストレージ装置とで前記識別子空間を分担する場合における相対的な空間幅を決定する第1の決定機能と、前記第1の決定機能により決定された前記相対的な空間幅に基づいて、前記識別子空間における前記第1の記憶機能に記憶された基準位置から前記隣接ストレージ装置に係る基準位置までの範囲内について分担する分担範囲を決定する第2の決定機能とを実現させるためのプログラムである。

#### 【0022】

なお、装置に係る本発明は方法に係る発明としても成立し、方法に係る本発明は装置に係る発明としても成立する。

また、装置または方法に係る本発明は、コンピュータに当該発明に相当する手順を実行させるための（あるいはコンピュータを当該発明に相当する手段として機能させるための、あるいはコンピュータに当該発明に相当する機能を実現させるための）プログラムとしても成立し、該プログラムを記録したコンピュータ読取り可能な記録媒体としても成立する。

#### 【0023】

本発明によれば、重みづけによって記憶容量の大きいノードや処理速度の速いノードにより多くの識別子空間を担当させ、より多くのデータを割り付けることによって、容量の小さいノードや処理の遅いノードがボトルネックになることを回避し、ノードの記憶容量や処理速度を有効に活用できる。

#### 【0024】

また、仮想化との組み合わせによってその割り付けはさらに正確に重みに比例させることができる。

#### 【0025】

また、空間の重複割り当てによる多重化、冗長化によって、ストレージノードの故障時には（すでにデータのレプリカを持っている）隣接ノードがただちに故障ノードの空間を継承して対応することができるので、高信頼化を実現できる。

#### 【0026】

また、本発明によれば、システム全体を集中制御するストレージノードなしに

、重みに比例等した空間割り当てを行なう分散アルゴリズムが実現でき、ストレージの使用率や応答速度が向上できる。

#### 【0027】

##### 【発明の実施の形態】

以下、図面を参照しながら発明の実施の形態を説明する。

##### （第1の実施形態）

図1に、本実施形態のシステム全体の構成例を示す。

図1において、1はストレージノード、3はクライアント計算機、7はネットワークを示す。なお、図1では、分散ストレージシステムを構成するストレージノードとして3台のものが示されているが、これは一例であり、分散ストレージシステムを構成するストレージノードの台数は任意である。また、図1では、分散ストレージシステム（ストレージノード群）が1グループだけ示されているが、複数グループ存在して構わない。この場合に、各ストレージノードはいずれか1つのグループにのみ属するようにする構成と、複数のグループに属するストレージノードが存在しても構わないものとする構成とがあり得る。

また、図1では、1台のクライアント計算機のみが示されているが、もちろん、クライアント計算機は複数台存在し得る。

#### 【0028】

さて、本実施形態の分散ストレージシステムは、各ストレージノードが重みを持ち、ファイルをストレージノードに対応付けるハッシュ空間の割り当ての際に、隣接するストレージノード間で、それらの重みに比例した幅にハッシュ空間を分割する。

#### 【0029】

なお、ハッシュ空間に両端が存在する構成の場合には、ハッシュ空間の端部を担当する2つのストレージノードはそれぞれ隣接するストレージノードを1つのみ持ち、それ以外の各ストレージノードはそれぞれ隣接ストレージノードを2つ持つことになる。また、ハッシュ空間の両端が連結されているものとみなしてハッシュ空間をループ状に形成する構成の場合には、いずれのストレージノードも隣接するストレージノードを2つ持つことになる。



## 【0030】

ハッシュ関数としては、ビット数が充分長くて一様な割り当てが行なわれる関数を用いると好ましい。そのようなハッシュ関数として、例えば、SHA-1を採用することができる。ハッシュ空間は、 $0 \sim 2^b - 1$ の $b$ ビットの整数で表わされる（SHA-1の場合は、 $b = 160$ となる）。

## 【0031】

本実施形態では、ネットワーク7上に分散配置された複数のストレージノード1が1つのストレージノード群（分散ストレージシステム）を構成する。前述のように、ある1つのストレージノード1が複数のストレージノード群に属する構成も可能であるが、その場合には、ハッシュ空間は、各ストレージノード群ごとに独立して管理するものとすればよい。

## 【0032】

なお、各ストレージノード1はストレージノード間通信に用いるアドレスを持つが、各ストレージノード1のノードID（ノード識別子）は、当該ストレージノード1の持つアドレスに対してハッシュ関数を適用して得られる値とする。

また、以下では、アドレスの一例として、IPアドレスを用いる場合を例にとって説明する。

## 【0033】

あるストレージノード群に属するストレージノード1は、それぞれ、自ノードの持つアドレスの他に、そのストレージノード群に属する他の全部のストレージノード1のアドレスを収集し、記憶するようにしてもよい。

## 【0034】

しかし、各ストレージノード1に他のすべてのストレージノード1のアドレスを記憶させると、例えばストレージノード1の増減の際に通信量および処理時間の問題が発生するような場合があるので、自ノードの持つアドレスの他に、そのストレージノード群に属する他の一部の（予め定められた条件を満たす）ストレージノード1のアドレスのみを収集し、記憶するようにしてもよい。例えば、各ストレージノード1には、ハッシュ空間上で自ノードに隣接するストレージノード1のアドレスと、ハッシュ空間上のいくつかの点（例えば、自ノードのノード

ID (自ノードのアドレスに対してハッシュ関数を適用して得られる値) に対して、1, 2, 4, 8, ...,  $2^{(b-1)}$  をそれぞれ加えた複数の点) を含むハッシュ空間を担当するノードのアドレスのみを記憶しておくようにしてもよい。これらのアドレスを利用すれば、 $O(b)$  個のノードへの問い合わせを行なうことによって、ハッシュ空間上の任意の点を担当するストレージノード 1 を検索することができる。

#### 【0035】

図 2 に、本実施形態のストレージノード 1 の内部構成例を示す。

図 2 に示されるように、ストレージノード 1 は、空間幅決定部 11、空間割当制御部 12、空間割当情報記憶部 13、ファイル入出力部 14、ファイル記憶部 15 を備えている。概略的には、空間幅決定部 11 と空間割当制御部 12 と空間割当情報記憶部 13 で、空間割り当ての決定や管理を行ない、ファイル入出力部 14 で、その空間割り当てに基づくファイルの入出力を行なう。また、ファイル記憶部 15 は、自ノードが分担するハッシュ空間に対応するファイルを記憶する。また、ストレージノード 1 は、自ノードに関する情報 (例えば、自ノードのアドレスやノード ID や重みなど) を記憶するための記憶部 (図示せず) を持つ。

#### 【0036】

図 3 に、本実施形態に係るストレージノードの参加時における処理手順の一例を示す。

ストレージノード 1 が、新規に、あるストレージノード群へ参加する場合は、まず、自ノードの IP アドレスを取得する (ステップ S1)。IP アドレスの取得方法には通常の IP ネットワークでの取得方法が利用できる。例えば、管理者がそのサブネット上で未使用となっている IP アドレス一つを適宜選択して手動で当該参加するストレージノード (以下、新ストレージノード) 1 へ入力するようにしてもよいし、あるいは、プールされた IP アドレスから一つを DHCP サーバが自動的に選択して DHCP プロトコルによって当該新ストレージノード 1 に通知するようにすることも可能である。

#### 【0037】

また、前述したように、当該新ストレージノード 1 は、入力あるいは通知され

た自ノードのアドレスに対してハッシュ関数を適用することによって、自ノードのノードIDを求める（ステップS2）。

#### 【0038】

次に、新ストレージノード1は、ノード群の任意の他の（例えば一つの）ストレージノードのIPアドレスを取得して、これを空間割当制御部12に伝える。他のストレージノードのアドレスの取得には、どのような方法を用いても構わない。例えば、管理者が他の（例えば一つの）ストレージノードを適宜選択してそのIPアドレスを手動で当該新ストレージノード1へ入力するようにしてもよいし、あるいは、近隣にストレージノードが存在する場合にDHCPのオプションやブロードキャストによって該近隣のストレージノードから当該新ストレージノード1へ、該近隣のストレージノードのIPアドレスを自動的に通知するようにすることも可能である。

#### 【0039】

このようにして得た他のストレージノードのIPアドレスを用いて、新ストレージノード1の空間割当制御部12は、適宜、他のストレージノード1（の空間割当制御部12）に接続して、自ノードの属するストレージノード群に属する他の全部又は一部のストレージノード1のアドレスを問い合わせ、収集する（ステップS3）。一部のストレージノード1のアドレスのみ記憶する場合には、例えば、自ノードのノードIDをもとに、隣接するノードのアドレスと、自ノードのノードIDに1, 2, 4, 8, ...,  $2^{(b-1)}$ をそれぞれ加えた点を含むハッシュ空間を担当する各ストレージノード1のアドレスを得る。

#### 【0040】

このようにして収集された他のストレージノード1のアドレスは、空間割当情報記憶部13に記憶される（ステップS4）。

#### 【0041】

以下、図4を参照しながら、本実施形態の各ストレージノードにおけるハッシュ空間分割方法について説明する。

以下の処理は、隣接ノードが2つある場合には、各々の隣接ノードについて行う。

## 【0042】

また、以下の処理を最初に行う場合には、例えば、図3のステップS4に続いて行ってもよいし、最初にクライアント計算機からI/O要求を受信した時点で行ってもよいし、それ以外の適当なタイミングで行ってもよい。なお、以下の処理は、例えば、自ノードの重みを変更された場合、隣接ノードの重みを変更された場合、ストレージノードの増減などによって隣接ノードが変更された場合などに、あらためて行うこともある。

## 【0043】

本実施形態のストレージノード1の空間幅決定部11は、空間割当制御部12から得た隣接ノードのアドレスをもとに、隣接するストレージノード1（の空間幅決定部11）との接続を行ない、自ノードの重みを伝えとともに、隣接するストレージノード1の重みを取得する。

## 【0044】

重みとしては、例えば、ストレージの記憶容量、演算能力、回線速度、それらを適宜組み合わせたものなどを用いることができる。ここでは、一例として、ストレージの記憶容量を用いた場合を例にとって説明する。

## 【0045】

自ノードsの重みをV[s]、隣接ノードuの重みをV[u]とすると、自ノードsの相対的な幅は、

$$w = V[s] / (V[s] + V[u])$$

となる。

## 【0046】

空間幅決定部11は、このようにして空間幅情報wを決定すると、これを、空間割当制御部12に送る。

空間割当制御部12では、空間幅情報wをもとに、自ノードsと隣接ノードuとの間の空間を分割する。なお、自ノードsと隣接ノードuとの間の空間は予め定められており、一方の端点が自ノードuのアドレスにハッシュ関数を適用して得られる点であり、他方の端点が隣接ノードsのアドレスにハッシュ関数を適用して得られる点である。

## 【0047】

しかして、例えば、ハッシュ関数を  $h()$  、自ノード  $s$  のアドレスを  $A[s]$  、隣接ノード  $u$  のアドレスを  $A[u]$  とすると、隣接ノード  $u$  との境界を

$$h1 = h(A[u]) * w + h(A[s]) * (1 - w)$$

とすることによって、

$$h1 - h(A[u]) : h(A[s]) - h1 = V[u] : V[s]$$

という条件を満たす境界  $h1$  が得られる。

この様子を、図4に示している。

## 【0048】

もう一方の隣接ノード  $d$  との境界  $h2$  についても、同様の手順によって、求めることができる。すなわち、もう一方の隣接ノード  $d$  のアドレスを  $A[d]$  とし、その重みを  $V[d]$  とすると、相対的な幅は、

$$w' = V[s] / (V[s] + V[d])$$

となり、隣接ノード  $d$  との境界を

$$h2 = h(A[d]) * w' + h(A[s]) * (1 - w')$$

とすることによって、

$$h2 - h(A[d]) : h(A[s]) - h2 = V[d] : V[s]$$

という条件を満たす境界  $h2$  が得られる。

## 【0049】

このようにして各隣接ノード  $u$  ,  $d$  との境界を決定し、その範囲 (図4の  $h1$  と  $h2$  参照) を、自ノード  $s$  の担当空間として、空間割当記憶部13に記憶する。

## 【0050】

なお、ハッシュ空間に両端が存在する構成の場合に、ハッシュ空間の端部を担当する2つのストレージノードについては、隣接するストレージノードとの境界を上記と同様の方法で求め、この境界から、ハッシュ空間の (自ノード側の) 端部までを、自ノードの担当空間とすればよい。

## 【0051】

このように重み (例えば、記憶容量) に応じて隣接ノード間でハッシュ空間を

分割する（例えば、重みに比例した幅で分割する）ことによって、図5に例示するように、例えば、記憶容量の大きなストレージノードにはより多くのファイルが割り当てられるようにすることが可能になる。

#### 【0052】

なお、上記では、空間幅決定部11で $w = V[s] / (V[s] + V[u])$ を求め、次いで、空間割当制御部12で $h1 = h(A[u]) * w + h(A[s]) * (1 - w)$ を求めるものとしたが、もちろん、 $h1 = h(A[u]) * V[s] / (V[s] + V[u]) + h(A[s]) * (1 - V[s] / (V[s] + V[u]))$ として求めるようにしてもよい。

#### 【0053】

また、上記した空間幅決定部11及び空間割当制御部12による決定方法は、一例であり、他にも種々の方法が可能である。

#### 【0054】

図6に、本実施形態に係るストレージノードの要求受信時における処理手順の一例を示す。

クライアント計算機1が分散ストレージシステムへファイルの読み書きを行う際には、クライアント計算機1は、自身が利用するファイルを管理しているノード群のストレージノード1に対してI/O要求を送出する。

#### 【0055】

なお、ノード群の任意のストレージノード1を利用することが可能であるが、実用上は、例えば、ネットワーク上で近い位置にあるストレージノード1を利用することが想定される。

#### 【0056】

ストレージノード1がI/O要求を受信すると（ステップS11）、そのファイル入出力部14では、空間割当情報記憶部13の自ノードの担当範囲と、要求中のファイルID（ファイル識別子）（なお、ファイルIDには、対象ファイルのファイル名あるいは対象ファイルの中身のデータに対して、ハッシュ関数を適用して得られる値などが用いられる）とを比較し（ステップS12）、ファイルIDが担当範囲内に該当すれば（ステップS13）、自ノードのファイル記憶部

15にアクセスして要求の処理（ファイルの読み出しあるいは書き込みなど）を行なう（ステップS14）。

#### 【0057】

他方、範囲外に該当すれば（ステップS13）、空間割当情報記憶部13の他のノードのアドレスをもとに、他のノード（の空間割当制御部12）に問い合わせることによって、そのファイルIDを担当するストレージノード1のアドレスを検索する（ステップS15）。

#### 【0058】

そして、検索によって、そのファイルIDを担当するストレージノード1のアドレスが得られたならば（ステップS16）、そのストレージノード1（のファイル入出力部14）に接続して、要求の処理を行なう（ステップS17）。

#### 【0059】

もし、そのファイルIDを担当するストレージノード1のアドレスが得られなかったならば（ステップS16）、エラー処理を行う（例えば、要求元のクライアント計算機に、エラー・メッセージを返す）（ステップS18）。

#### 【0060】

このように本実施形態によれば、空間幅を重みに応じたものにするによって、例えば、ストレージノードの記憶容量などに応じたファイルの割り当てや、処理速度に応じた負荷の分散などが可能になる。

#### 【0061】

（第2の実施形態）

本実施形態の分散ストレージシステムは、1つのノード群を構成するストレージノード（の全部又は一部）を複数の仮想ノードに仮想化し、複数のハッシュ空間に割り当てた上で、ファイルへの対応付けを行なうものである。

#### 【0062】

以下では、第1の実施形態と相違する点を中心に説明する。

#### 【0063】

図7に、本実施形態のハッシュ空間分割方法を示す。

#### 【0064】

各ストレージノード（以下、仮想ノードと区別するため、物理ノードと呼ぶ）  
1 のノード ID からハッシュ値を計算する際に、物理ノードのアドレスに加えて、物理ノードごとに割り当てられる仮想ノード番号をハッシュ関数の引数として用いる。

#### 【0065】

例えば、ある物理ノード  $s$  のアドレスが  $A[s]$ 、物理ノード  $s$  に設定される仮想ノードの数が  $v$  個である場合に、物理ノード  $s$  には、 $h(A[s], 0)$ ,  $\dots$ ,  $h(A[s], v-1)$  の  $v$  個のハッシュ値が対応する。

#### 【0066】

そして、各仮想ノードをそれぞれ第 1 の実施形態のストレージノードとみなせば、第 1 の実施形態と同様の処理によって、各仮想ノードのハッシュ空間の担当範囲を求めることができる。

#### 【0067】

なお、仮想ノードごとに重みを付与する方法と、全仮想ノードについて共通の重みを付与する方法が可能である。

#### 【0068】

本実施形態は、第 1 の実施形態の利点に加えて、物理ノードに複数の空間を分割させることによって、より多くのノードに隣接させることができるため、隣接ノードの重みの平均値が全ノードの重みの平均値に近づくようになり、したがって、各ノードの担当する空間幅がより正確に重みに比例するという利点がある。これは、重み付けを行わない従来技術に仮想化を適用した場合に得られる、ノード間のハッシュ空間上での距離の分散が小さくなるという効果とは異なり、重みに応じた空間分割と仮想化とを組合わせたからこそ得られる効果である。

#### 【0069】

（第 3 の実施形態）

本実施形態に係る分散ストレージシステムは、空間を複数のノードに重複して割り当て、ファイルをそれらのノードに重複して割り当てるものである。

#### 【0070】

以下では、第 1 の実施形態と相違する点を中心に説明する。



**【0071】**

図8に、本実施形態のハッシュ空間分割方法を示す。

**【0072】**

なお、図8は、 $n=2$ （二重化）の場合を例示している。

**【0073】**

ハッシュ空間の同一の点を  $n$  個のノードに重複して割り当てる場合、ストレージノードを  $n$  グループ  $n$  分割し、同一のハッシュ空間を対象として、各々のグループ内で、第1の実施形態と同様に、各ストレージノードのハッシュ空間の担当範囲を求めればよい。

**【0074】**

その際、例えば、第1の実施形態の隣接するノードの代わりに、ハッシュ空間上で自ノードの  $n$  個隣のストレージノードとの間で、担当空間の分割を行なうようにしてもよい。それらに挟まれた他のストレージノードも同様の空間を分割するため、結果として、一つの点が  $n$  個ノードに重複して割り当てられることになる。

**【0075】**

ところで、単に多重化するだけであれば、ファイルを複数のノードに対応させる種々の方法がある。例えば、ファイルに複数のハッシュ関数を適用することによって、複数個の仮想的なファイルIDを割り当て、それらの仮想的なファイルIDに対応するノードIDを持つ複数個のノードにファイルを格納する、といった冗長化方法が考えられる。しかしながら、故障などによってストレージノードが失われた場合に、そのストレージノードが担当していた空間を引き継ぐのは、ハッシュ空間上で故障ノードに隣接するストレージノードである。このような方法では、故障ノードに隣接するストレージノードに、他のストレージノードから該当するファイルを転送しなければならない。

**【0076】**

これに対して、本実施形態では、隣接するストレージノードに重複した割り当てが行なわれているため、そのストレージノードに空間を継承させることによって、ファイルの転送量を最小化することができる。

**【0077】**

本実施形態によれば、第1の実施形態の利点に加えて、ファイルを複数のノードに格納することによって、故障などでファイルが失われる可能性を低減させるとともに、故障などの際に、ファイルの転送量を最小限に抑えて担当空間の継承ができるという利点がある。

**【0078】**

なお、第3の実施形態と第4の実施形態とを組み合わせることも可能である。

**【0079】**

なお、以上では、ファイルをストレージノードに格納するにあたって、ファイル単位で格納する場合を例にとって説明してきたが、ファイルをブロック単位に分割してブロック単位で格納する場合も同様に可能である。この場合には、例えば、ファイル名とブロック番号との組をブロック名とするなどによって、これまで説明してきたものと同様の方法が適用可能になる。

**【0080】**

なお、以上の各機能は、ソフトウェアとして実現可能である。

また、本実施形態は、コンピュータに所定の手段を実行させるための（あるいはコンピュータを所定の手段として機能させるための、あるいはコンピュータに所定の機能を実現させるための）プログラムとして実施することもでき、該プログラムを記録したコンピュータ読取り可能な記録媒体として実施することもできる。

**【0081】**

なお、この発明の実施の形態で例示した構成は一例であって、それ以外の構成を排除する趣旨のものではなく、例示した構成の一部を他のもので置き換えたり、例示した構成の一部を省いたり、例示した構成に別の機能あるいは要素を付加したり、それらを組み合わせたりすることなどによって得られる別の構成も可能である。また、例示した構成と論理的に等価な別の構成、例示した構成と論理的に等価な部分を含む別の構成、例示した構成の要部と論理的に等価な別の構成なども可能である。また、例示した構成と同一もしくは類似の目的を達成する別の

構成、例示した構成と同一もしくは類似の効果を奏する別の構成なども可能である。

また、この発明の実施の形態で例示した各種構成部分についての各種バリエーションは、適宜組み合わせて実施することが可能である。

また、この発明の実施の形態は、個別装置としての発明、関連を持つ2以上の装置についての発明、システム全体としての発明、個別装置内部の構成部分についての発明、またはそれらに対応する方法の発明等、種々の観点、段階、概念またはカテゴリに係る発明を包含・内在するものである。

従って、この発明の実施の形態に開示した内容からは、例示した構成に限定されることなく発明を抽出することができるものである。

#### 【0082】

本発明は、上述した実施の形態に限定されるものではなく、その技術的範囲において種々変形して実施することができる。

#### 【0083】

##### 【発明の効果】

本発明によれば、分散ストレージシステムを構成する各ストレージ装置が、該分散ストレージシステムが対象とする識別子空間のうちで分担する分担範囲を決定するにあたって、より効果的な分担範囲の決定を行うことができる。

##### 【図面の簡単な説明】

【図1】 本発明の一実施形態に係る分散ストレージシステムの全体構成例を示す図

【図2】 同実施形態に係るストレージノードの内部構成例を示す図

【図3】 同実施形態に係るストレージノードの参加時における処理手順の一例を示すフローチャート

【図4】 同実施形態のハッシュ空間分割方法を説明するための図

【図5】 同実施形態のハッシュ空間分割方法を説明するための図

【図6】 同実施形態に係るストレージノードの要求受信時における処理手順の一例を示すフローチャート

【図7】 同実施形態の他のハッシュ空間分割方法を説明するための図

【図 8】 同実施形態のさらに他のハッシュ空間分割方法を説明するための  
図

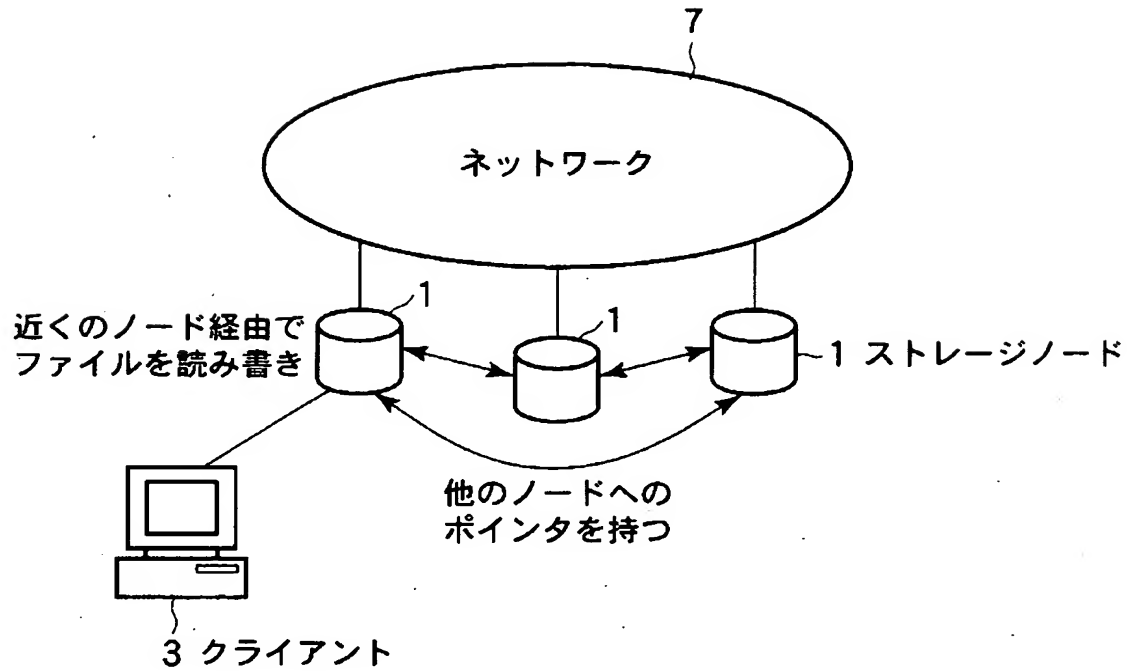
【符号の説明】

1…ストレージノード、3…クライアント計算機、7…ネットワーク、11…  
空間幅決定部、12…空間割当制御部、13…空間割当情報記憶部、14…ファ  
イル入出力部、15…ファイル記憶部

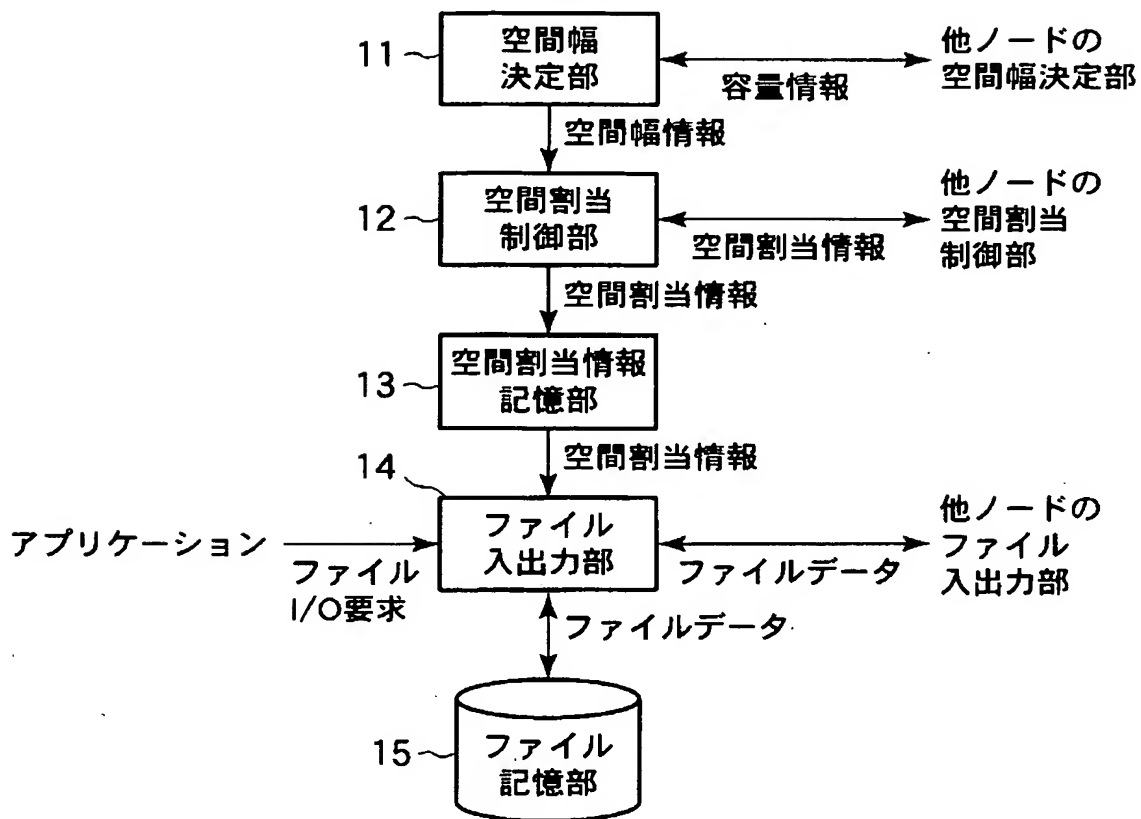
【書類名】

図面

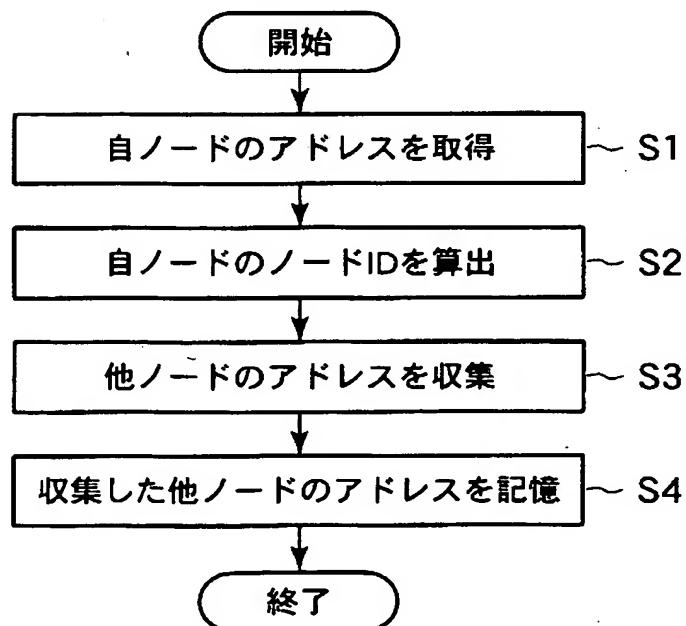
【図 1】



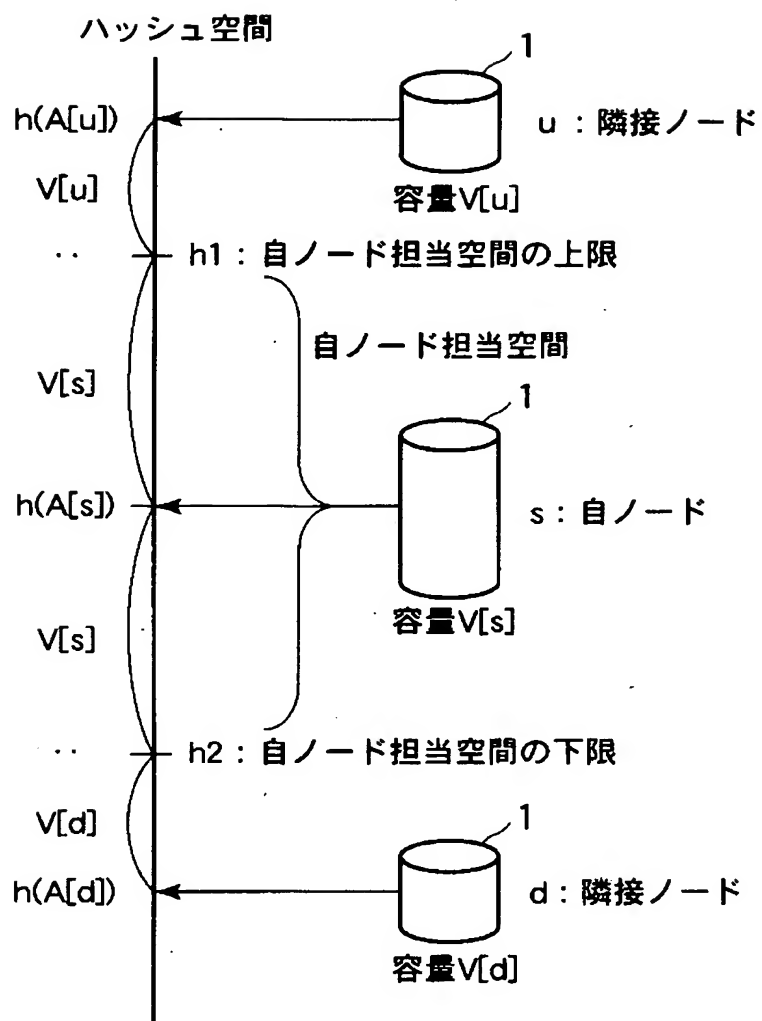
【図 2】



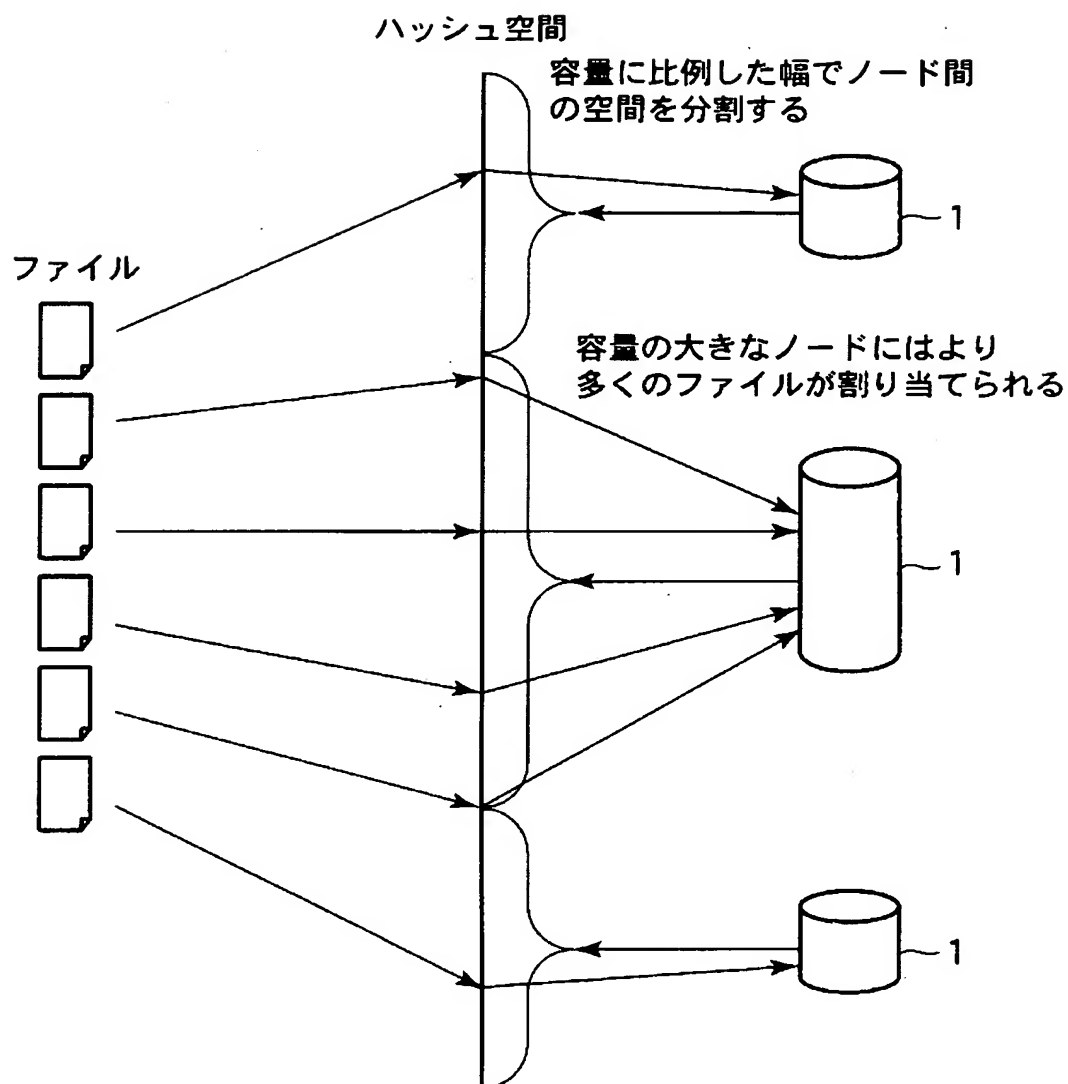
【図 3】



【図 4】

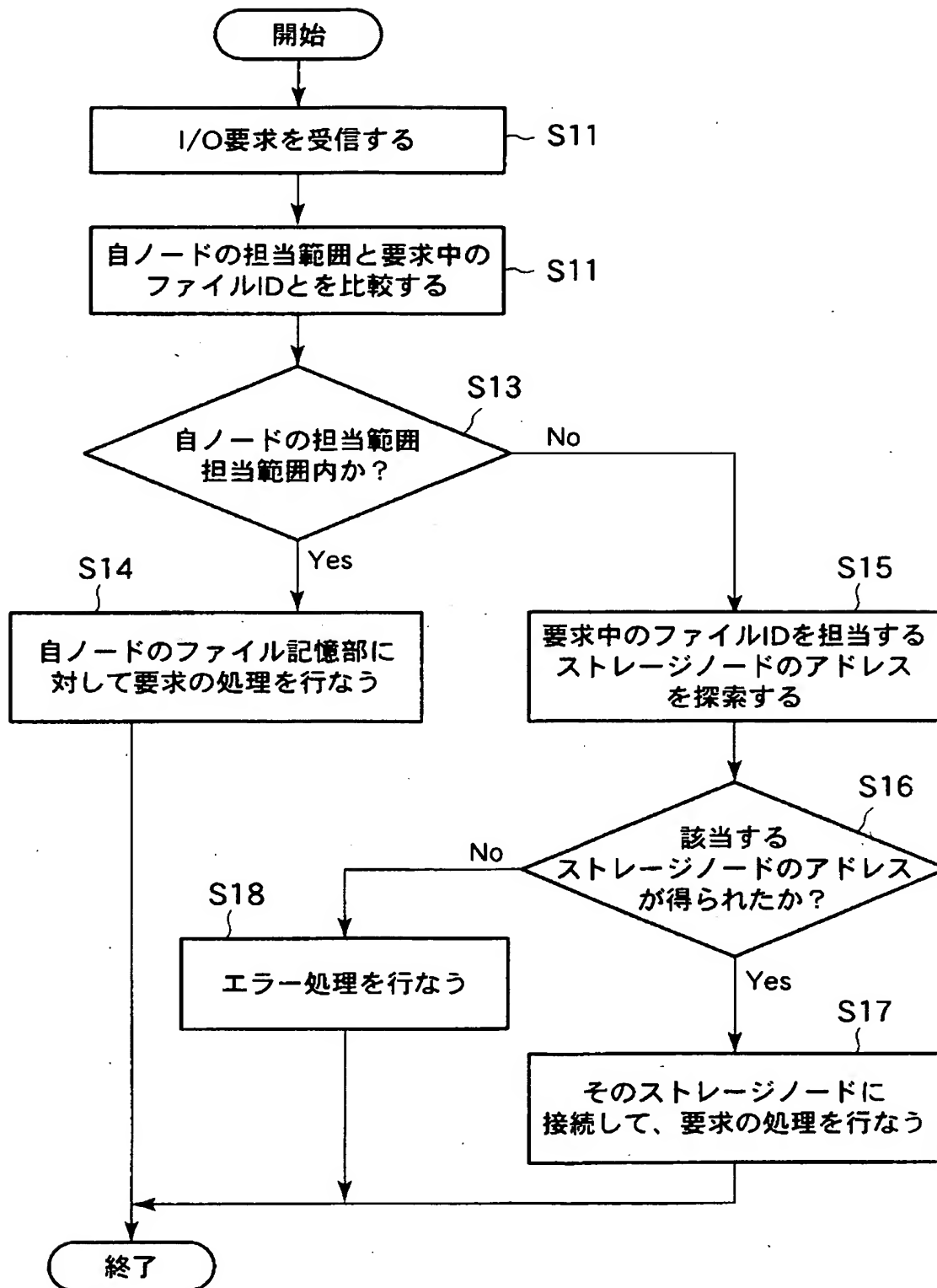


【図 5】

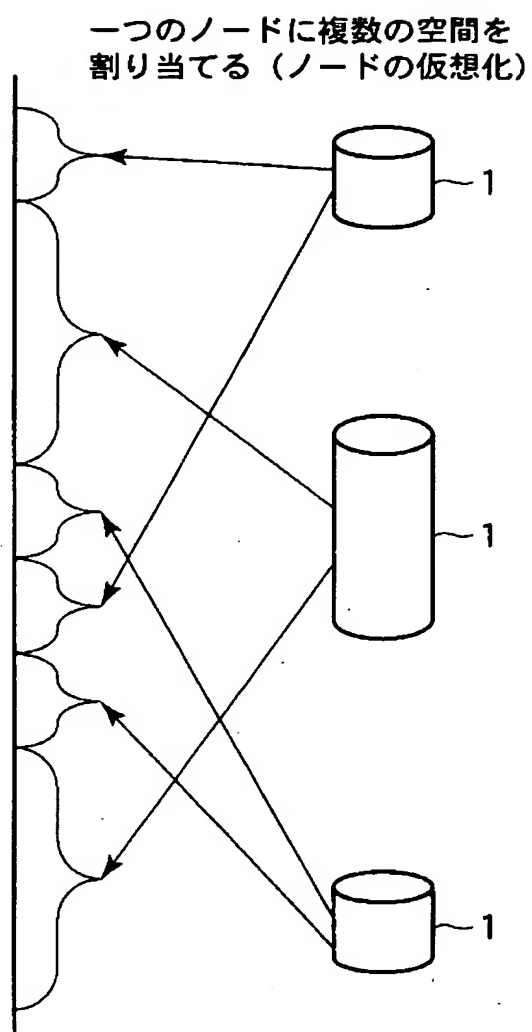




【図 6】

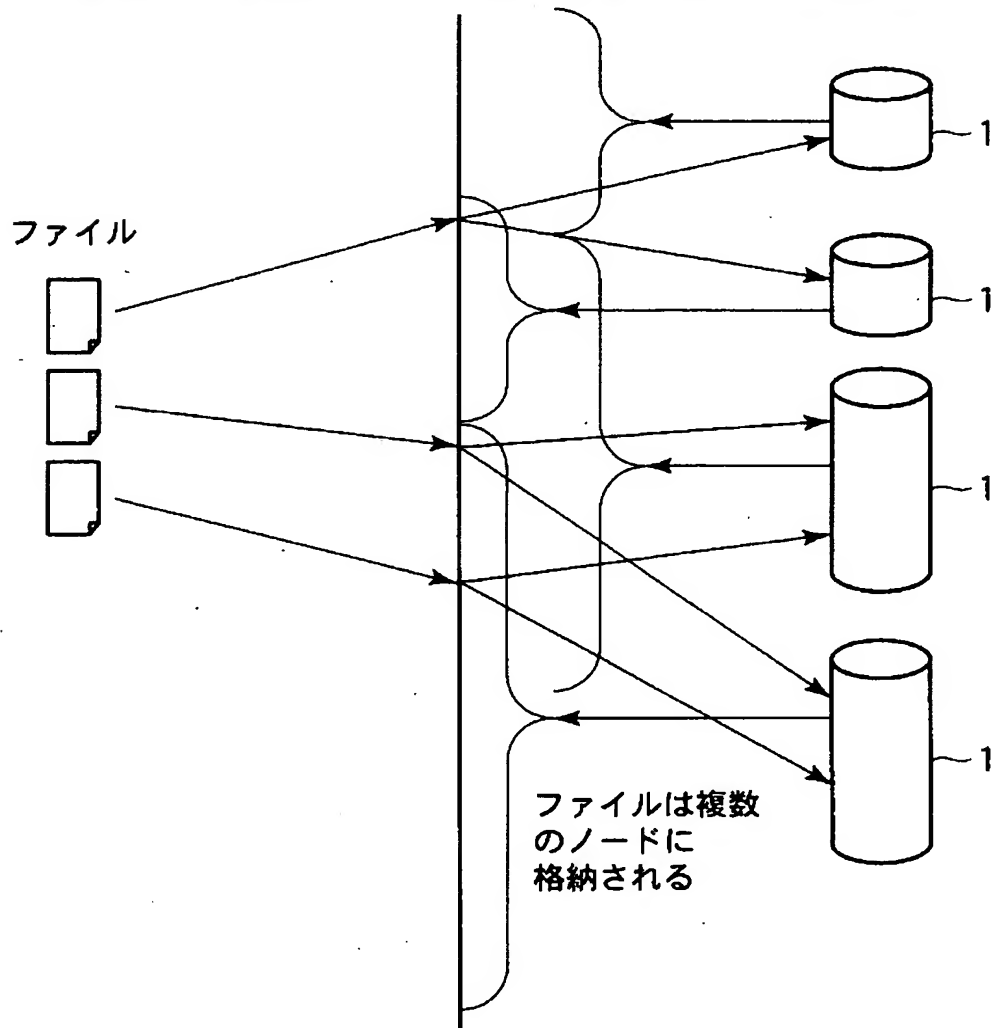


【図 7】



【図 8】

多重化のために空間を重複して割り当てる  
(二重化なら二個隣のノードとの間で容量に比例した空間分割)



【書類名】 要約書

【要約】

【課題】 分散ストレージを構成する各ストレージ装置が、分散ストレージが対象とするファイル識別子空間のうちで分担する範囲を決定するにあたって、より効果的な分担範囲の決定を行うことができるストレージ装置を提供すること。

【解決手段】 ストレージ装置のファイル記憶部 1 5 は、分散ストレージシステムが対象とするファイル識別子空間のうちで自ノードが分担する分担範囲に含まれるファイル識別子を持つファイルを格納する。空間幅決定部 1 1 は、自ノードの持つ重みと、隣接ノードの持つ重みとから、相対的な空間幅を決定する。空間割当制御部 1 2 は、自ノードのアドレスに所定のハッシュ関数を適用して、ファイル識別子空間における第 1 の基準位置を求め、同様に、隣接ノードのアドレスから第 2 の基準位置を求める、相対的な空間幅に基づいて、第 1 の基準位置から第 2 の基準位置までの範囲内について分担する分担範囲を決定する。

【選択図】 図 1

特願 2003-041486

出願人履歴情報

識別番号

[000003078]

- |          |                |
|----------|----------------|
| 1. 変更年月日 | 2001年 7月 2日    |
| [変更理由]   | 住所変更           |
| 住 所      | 東京都港区芝浦一丁目1番1号 |
| 氏 名      | 株式会社東芝         |
|          |                |
| 2. 変更年月日 | 2003年 5月 9日    |
| [変更理由]   | 名称変更           |
|          | 住所変更           |
| 住 所      | 東京都港区芝浦一丁目1番1号 |
| 氏 名      | 株式会社東芝         |